



# P4Visor: Lightweight Virtualization and Composition Primitives for Building and Testing Modular Programs

Peng Zheng  
Xi'an Jiaotong University and  
Brown University  
zeepean@gmail.com

Theophilus Benson  
Brown University  
theophilus\_benson@brown.edu

Chengchen Hu  
Xi'an Jiaotong University  
huc@ieee.org

## ABSTRACT

Programmable data planes, PDPs, enable an unprecedented level of flexibility and have emerged as a promising alternative to existing data planes. Despite the rapid development and prototyping cycles that PDPs promote, the existing PDP ecosystem lacks appropriate abstractions and algorithms to support these rapid testing and deployment life-cycles. In this paper, we propose P4Visor, a lightweight virtualization abstraction that provides testing primitives as a first-order citizen of the PDP ecosystem. P4Visor can efficiently support multiple PDP programs through a combination of compiler optimizations and program analysis-based algorithms. P4Visor's algorithm improves over state-of-the-art techniques by significantly reducing the resource overheads associated with embedding numerous versions of a PDP program into hardware. To demonstrate the efficiency and viability of P4Visor, we implemented and evaluated P4Visor on both a software switch and an FPGA-based hardware switch using fourteen different PDP programs. Our results demonstrate that P4Visor introduces minimal overheads (less than 1%) and is one order of magnitude more efficient than existing PDPs primitives for concurrently supporting multiple programs.

## CCS CONCEPTS

• **Networks** → **Programmable networks**; • **Software and its engineering** → *Software testing and debugging*;

## KEYWORDS

Programmable Data Plane, Code Merge, Testing

### ACM Reference Format:

Peng Zheng, Theophilus Benson, and Chengchen Hu. 2018. P4Visor: Lightweight Virtualization and Composition Primitives for Building and Testing Modular Programs. In *CoNEXT '18: International Conference on emerging Networking EXperiments and Technologies, December 4–7, 2018, Heraklion, Greece*. ACM, New York, NY, USA, 14 pages. <https://doi.org/10.1145/3281411.3281436>

## 1 INTRODUCTION

Programmable data planes [10, 39, 41] (PDPs), e.g., Tofino [39], have emerged as a promising alternative to traditional data planes. These

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from [permissions@acm.org](mailto:permissions@acm.org).  
*CoNEXT '18, December 4–7, 2018, Heraklion, Greece*

© 2018 Association for Computing Machinery.  
ACM ISBN 978-1-4503-6080-7/18/12...\$15.00  
<https://doi.org/10.1145/3281411.3281436>

PDPs enable an unprecedented level of flexibility: they provide abstractions and language frameworks that simplify the development of stateful network functionality which operates at line rate. This flexibility enables rapid development and prototyping of novel functionality and use cases.

Despite these rapid development and prototyping cycles, the existing PDP ecosystem lacks appropriate primitives and algorithms to support rapid testing and deployment. At a high level, many testing paradigms [31, 52, 59], e.g., canary testing used in Google's [20, 46] networks, require running new versions of a program alongside stable versions. Traffic is split across all versions and the output is compared. Orthogonally, supporting agile development requires composing and merging modular programs together.

The key challenges to enabling these techniques in today's PDP networks lie in efficiently supporting multiple PDP programs and providing flexible operators for the broad range of potential paradigms. Hardware PDP devices include limited physical resources which restrict the size of the PDP programs that can be supported, and enabling multiple versions of a PDP programs on a resource constraint device requires effective algorithms for minimizing resource footprints. Additionally, PDP language abstractions, e.g., P4, provide a limited set of primitives, e.g., P4 does not support loops, and the language restrictions complicate the process of developing general primitives to support a broad range of scenarios. Specifically, in this paper, we focus on one of the most popular and promising data plane programming languages – P4.<sup>1</sup>

In this paper, we present P4Visor, an abstraction layer and composition primitives, which addresses the above challenges to make testing and development primitives first-order citizens of the PDP ecosystem. The key insight behind P4Visor is that the different versions of a P4 program will share significant code fragments (i.e., tables, parse graph states and action primitives) and thus we can reduce the resource overheads by merging the P4 programs and thus eliminating redundancy. In this way, an administrator can run multiple P4 programs concurrently in the data plane.

P4Visor achieves this through a combination of program analysis to identify potential program overlaps and compiler optimizations to merge the P4 programs and reduce resource footprints. To flexibly support different testing paradigms, P4Visor includes domain-specific comparator operators that provide building blocks for composing new testing paradigms.

Today, the prevalent approach for supporting multiple P4 programs is to virtualize the data plane [22, 56], e.g., Hyper4 [22], HyperV [56], and host different programs atop the virtualization layer. Unfortunately, these approaches [22, 56] require significant

<sup>1</sup> The PDP Programs in the following sections refer to P4 programs unless otherwise stated.

resources and are often slow or unscalable [22] because they provide *Full-Virtualization* which uses software to emulate hardware. Rather than providing virtualization and modularity primitives, at the software layer, we aim to provide these primitives at the compilation layer which enables us to explore tradeoffs between flexibility and efficiency. In particular, our design choices allow us to trade off a modest amount of flexibility for a significant increase in efficiency.

Logically, P4Visor operates between the PDP programs and the PDP hardware providing resource management between different PDP programs and merge capabilities to the PDP Programs. It provides virtualization primitives required for supporting concurrent testing in production networks. P4Visor’s goals include security isolation between the management functionality provided by P4Visor’s interfaces and data plane functions running on the PDPs devices; efficient resource utilization and management; and, flexible support over arbitrary PDPs targets.

We make the following contributions:

- **Virtualization Abstractions:** We provide an abstraction for seamlessly merging two programs to tackle the resource management and indirection challenges that arise from merging and composing programs (Section 3).
- **Composition Operators:** We introduce several composition operators for merging P4 programs to support a range of testing paradigms (Section 3.2).
- **Merge Algorithm:** We present a first look at the code-merging problem for P4 programs. We build a model to theoretically identify the key issues behind merging P4 programs and demonstrate the hardness of the problem – it is NP-Complete. As a result, we propose a heuristic to solve it effectively (Section 4 & Section 5).
- **Prototype Implementation and Evaluation:** We implement a prototype of P4Visor’s framework and merging algorithms. Using this prototype, we demonstrate the flexibility and efficiency of P4Visor by testing it across multiple P4 Programs (Section 6 & Section 7).

## 2 MOTIVATION

In this section, we describe several well-understood principles used within production networks (*e.g.*, Google and Facebook) to ensure highly-available networks (Section 2.1), and present a new PDP primitive for effectively supporting these principles in PDP devices (Section 2.2 and Section 2.3).

### 2.1 Rapid Development in Large Networks

We briefly describe several techniques which large-scale networking infrastructures employ to ensure that their networks remain highly-available in the face of changes.

**Canary Testing (A-B Testing [3, 52, 59])** Canary testing, well documented by Google’s [20, 46, 59] and Facebook’s [52] networking and infrastructure teams, requires running multiple versions of a program alongside each other. Canarying (or A-B Testing) tests new code by sending a subset of traffic through the code (*e.g.*, 1% of traffic) and, if nothing “bad” happens, slowly increases the subset of traffic using the test code until all traffic is using the test code.

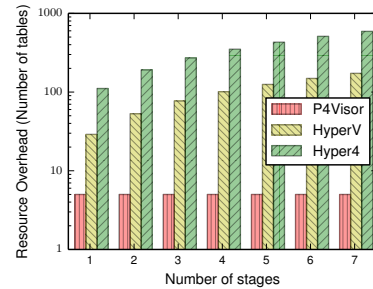


Figure 1: Comparison of virtualization overheads.

**Fault Tolerance (Data-Diversity [25, 31])** To improve security and availability, certain networks run multiple instances of their control plane, perturb the instances with some randomness, and then compare the outputs from these versions. The system uses the most popular output. This approach directly tackles bugs and overcomes intruders. Facebook [25] runs four control planes and compares the output between these control planes.

**Modular Code** Extensive work in the software engineering community [42] and recently validated by large software engineering firms (*e.g.*, Facebook [51, 52]) have demonstrated that the key to successfully supporting rapid prototyping and deployment of complex functionality is modularity (code-reuse). Yet, today programmers are forced to write monolithic P4 programs. Missing from the ecosystem is a framework for effectively supporting multiple modular PDP Programs – similar to processes – and composing them together.

### 2.2 Novel PDP Primitives: Code Merge

Today, the most direct approach for supporting the aforementioned testing techniques is to use virtualization, *e.g.*, HyperV and Hyper4. Unfortunately, HyperV and Hyper4 incur significant performance and resource overheads. In Figure 1, we present the memory overheads of using these different virtualization techniques with an emphasis on the number of tables used. The overheads grow linearly with the size of the program because both techniques declare a fixed number of additional tables to emulate each of the P4 program’s stages and primitive actions – their hypervisors have to use these tables to record the runtime states for each program. For example, to run a P4 program with two pipeline stages, HyperV and Hyper4 have to declare at least 53 and 191 tables respectively which limits the number of primitive actions supported to 9 and prevents HyperV and Hyper4 from supporting Switch.P4 [14] which has 19 primitive actions.

Motivated by the inefficiencies of existing virtualization primitives, in this paper we aim to answer the following question.

*Is it possible to have a framework for supporting multiple versions without incurring the overheads of full virtualization?*

To answer this question, we investigate the design of a light-weight virtualization based on a source code merging primitive. The merge primitive takes as input  $N$  P4 programs and creates as output one P4 program that combines all input P4 programs but retains the functionality of each of the original P4 programs. As an example shown in Figure 3, our new primitive takes, as input,

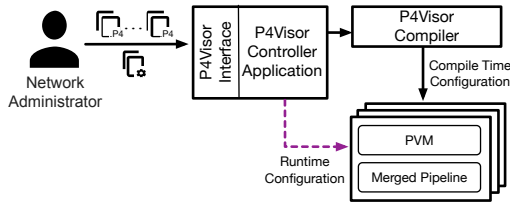


Figure 2: P4Visor workflow.

the abstract representations of two P4 programs (Production and Testing P4 control programs in Figure 3) and combines them into one (Figure 3 (a)). Central to providing this primitive is ensuring that during the merge, P4-specific correctness constraints (e.g., table dependencies) are maintained while efficiency is maximized through resource reuse.

The merged P4 program should give each P4 program the illusion of sole occupancy on the hardware. Our approach differs from full virtualization in several ways: first, while full virtualization provides virtualization through a special P4 program, we provide virtualization through the P4 program compiler. Our approach offers one key benefit: whereas full virtualization needs to allocate resources to support any potential P4 program, we only need to allocate resources to support the P4 programs being compiled. This specialization minimizes the number of additional tables required to support the combined program. Second, full virtualization does not explicitly support modularity and composition of multiple P4 programs into one, whereas, we can directly support these use cases.

### 2.3 P4Visor Workflow

Next, we present the workflow of P4Visor to illustrate P4Visor’s components and how they work together. To support the envisioned testing paradigms, the network operators must provide P4Visor with (1) the different P4 programs to merge, (2) the type of testing composition operators to implement (e.g., A-B testing or Differential testing) – the composition operator determines the policy for splitting and comparing traffic, (3) the amount of traffic used for testing, e.g., test X% of the traffic, and (4) the traffic sampling granularity, e.g., all packets of a flow should be consistently tested, test any packet of any flow, or just test flows within a specific subnet. Network administrators configure these settings using either a simple command line or a configuration file. Given such a testing specification, P4Visor installs code at edge switches to consistently tag packets for testing and to remove the tags before the packets exit the network. Tagging at the edge enables P4Visor to ensure that the different switches consistently test the same packets and that we can perform end-to-end tests across the whole network. During the merge, P4Visor adds tables to compare the output from the different program and generate packets to the controller when these results differ. These packets allow operators to reason about the implications of the new code. In Figure 2, we present P4Visor’s workflow.

To support this workflow, P4Visor requires (1) a domain-specific configuration language for configuring the testing paradigms (the P4Visor interface), (2) a merge algorithm (discussed in Section 4 and

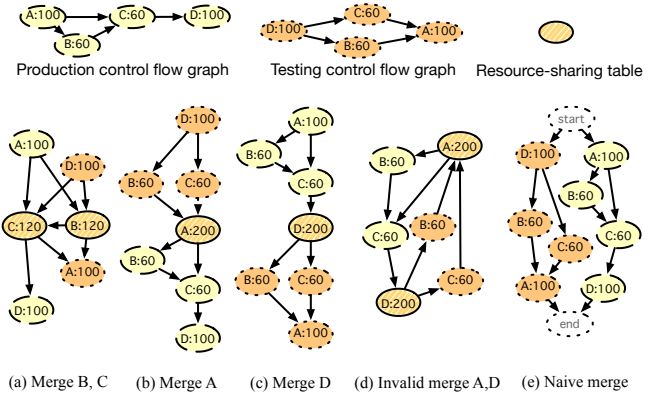


Figure 3: Illustrates various approaches for merging of two P4 programs: (a) demonstrates an intelligent merge with two share tables; (b) and (c) are two possible merges with one shared table; (d) is an invalid merge; and (e) demonstrates a simple combination of two programs which doubles the resources.

Section 5), (3) a framework for implementing the merge and supporting the indirection required to support the merge (discussed in Section 3), (4) techniques for enabling comparisons and techniques for tagging/untagging the packets.

## 3 DESIGN OF P4VISOR

In this section, we provide an overview of P4Visor’s architecture (Section 3.1), the currently supported composition operators (Section 3.2 & 3.3), and P4Visor’s compiler design (Section 3.4).

### 3.1 Overview

In Figure 4, we present the architecture for P4Visor. At a high-level, P4Visor is composed of four components: the P4Visor Compiler (PVC), P4Visor Management agent (PVM), the P4Visor controller Application (PVA), and P4Visor Interface (PVI).

**P4Visor Interface (PVI):** PVI runs on a server and provides the management interface for the network administrator (or developers) to use to control the composition of different P4 programs. In our current prototype, we implemented two operators: *A-B Testing* and *Differential Testing* (described in Section 3.2). As part of future work, we will explore other merge operators that are more amenable to modularization (e.g., Parallel composition [38], Sequential composition [38] ) and availability (e.g., Data-deduplication [31]).

**P4Visor Compiler (PVC):** PVC takes, as input the P4 programs, from the PVI, and returns, as output, a merged P4 program. The PVC analyzes the parse graphs, tables and control flows of the input P4 programs and merges them. The key to the PVC is the P4 program-merge algorithm (Section 4 & Section 5), which identifies the data plane resources within all input P4 programs that can be “safely” merged while maintaining the semantics and dependencies of each P4 program.

As a result of the merge process, P4Visor creates a new P4 program, which is a normal P4 program that can be run through a standard PDP compiler. Additionally, P4Visor creates a P4Visor-specific file, called the P4VisorConfiguration, which provides a

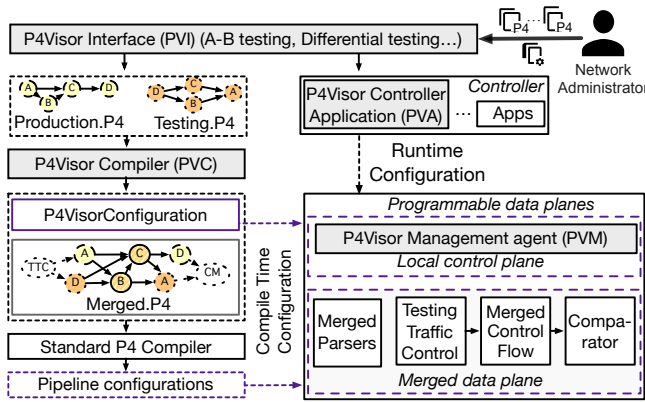


Figure 4: P4Visor overview.

mapping between the resources of each input P4 program and the merged P4 programs. Recall, each P4 program uses unique IDs to identify tables, and during the merge, these IDs may be modified. The P4VisorConfiguration provides a mapping between the original ID and the modified IDs: this file allows P4Visor to transparently rewrite all calls between the control and data plane that use these IDs. For example, when the controller sends a flow entry installation message to the switch to add a new entry to a *resource-sharing* table, PVM will modify the table IDs according to the P4VisorConfiguration to ensure that the entry is installed correctly.

**P4Visor Management agent (PVM):** PVM runs on the PDP devices, *i.e.*, programmable switch, intercepts messages between the control plane (controllers) and the merged P4 Program and uses the P4VisorConfiguration to determine how to appropriately modify the messages. Essentially, the agent multiplexes and demultiplexes messages between the different controllers and the merged P4 program.

**P4Visor controller Application (PVA):** PVA runs on the controller with a global view of the network, providing runtime control over the testing operators. For example, in A-B Testing, the PVA populates the testing traffic control tables (Section 3.3) for all the edge switches, identifies testing traffic.

### 3.2 Composition Operators

To illustrate the flexibility of our merge algorithm, we use it to implement two distinct testing composition operators.

**A-B Testing Operator:** This operator allows multiple programs to run side by side in a production network with a subset of traffic siphoned to the testing version, as shown in Figure 5 (a). To support, our A-B Testing composition, P4Visor must securely and flexibly manage traffic among multiple versions. To ensure security, *i.e.*, production traffic will not be processed by the testing programs, P4Visor adds/removes a special flag (TFlag) to packets at edge switches when traffic enters and leaves the network.

**Differential Testing Operator:** The key difference between the A-B Testing and Differential Testing operators is that: while the A-B Testing operator is mutually exclusive, *i.e.*, traffic either

goes to the production or the test P4 program, for the Differential Testing operator, the test packets must be copied and sent through both programs, with outputs compared at the end of the pipeline. The packet life-cycle is shown in Figure 5 (b).

### 3.3 Primitives for Composition Operators

To support these two operators, P4Visor must provide a flexible primitive for controlling traffic and, specifically, for Differential Testing, P4Visor must provide primitives for performing comparisons.

**Flexible Control:** To ensure flexibility, a traffic management module, called Testing Traffic Control (TTC in Figure 5), is developed and inserted into the merged pipeline to identify a packet as either “test” or “production” traffic and guide the packets along the appropriate pipeline. As shown in Figure 5, the TTC module is instantiated within the first table that all packets encounter and affixes the TFlag header to the packet once the packet is identified as the testing packet.

The TTC contains a set of stateful registers and flow tables, which are configurable using the P4Visor Interface. Using the P4Visor interface, network operators can configure the TTC to configure how traffic is sampled for testing:

- *Random sampling:* Operators can specify which percentage of traffic is randomly selected and piped through the testing “pipeline” and which percentage of traffic goes through the production “pipeline”, *e.g.*, randomly sample 1% of the total traffic for testing.
- *Flow based sampling:* Alternatively, operators can specify the exact flows that should be sampled by specifying a flowspect, *e.g.*, traffic from subnet “10.10.10.0/24” should be sampled.

**The Comparison Primitive:** To enable comparisons, a special table, called an output record table, is added at the end of each program’s pipelines. This table’s fields are configurable through the P4Visor interface. Specifically, we need to clone the packets and, in turn, compare their outputs at each switch.

To clone packets, we leverage the recirculate primitive, which recirculates a packet and allows the packet to be processed multiple times by the switch. We recirculate once for each version we want to test – during recirculation, we also recirculate the metadata fields. At the end of the pipeline, the packet is processed by the Comparator Module (the CM in Figure 5), which compares the output of the different versions. The comparator reports to the controller, via a message, when the compared packets and metadata fields are different.<sup>2</sup>

Using the P4Visor Interface, a network administrator can 1) specify which outputs should be compared; 2) configure how to process the differences detected by the Comparator. In general, the Comparator can support two kinds of comparisons either on packet header fields or on metadata fields.

<sup>2</sup>To control the overheads associated with recirculating packets, the operators can fine-tune the number of packets sampled to ensure the system provides acceptable performance.



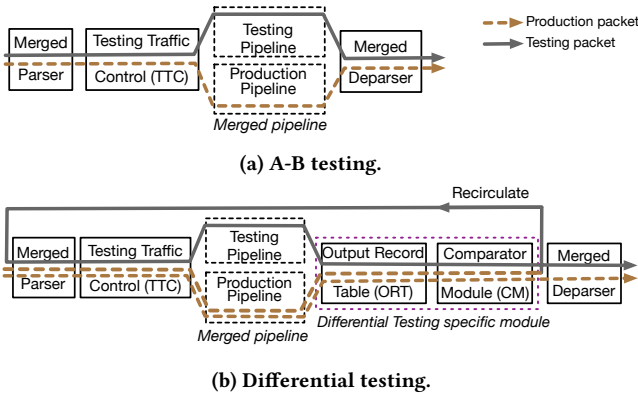


Figure 5: Life-cycle of a packet in P4Visor (under our two composition operators).

### 3.4 P4Visor Compiler

P4Visor merges PDP programs by merging their parse graphs and their control flow graphs:

**Merging Parse Graphs:** A naive approach for merging the parse graphs of two programs is to emulate the merge process by resubmitting packets into the pipeline multiple times; once for each of the states in both parse graphs (this is the approach taken by Hyper4 [22]). Unfortunately, this incurs a significant performance overhead – each time a packet is recirculated, throughput is cut and packet processing latency is increased.

Instead, we merge the two parse graphs into one and use a tag (*i.e.*, TFlag) to disambiguate and break conflicts in the merged parse graph. Figure 6 presents an example of two parse graphs being merged. Recall, each parse graph is a finite state machine (FSM) with each state representing the bit offsets of each header type. In merging these parse graphs, we align the parse graphs’ FSMs and merge identical states.<sup>3</sup> In Figure 6 (c), the merged states are in orange with solid line. We observe that Ethernet and IPv4 are both identical and thus can be merged. There is ambiguity about when to parse the VLAN and the IPv6 headers, and we break this ambiguity by introducing the TFlag state: packets with the TFlag, *i.e.*, test packets, should parse the IPv6 header type, whereas only packets without the TFlag, *i.e.*, non-test packets, should parse the VLAN header type. Note: we need only insert one such state for the TFlag, and this state can disambiguate all potential ambiguities in all merged states.

**Merging Control Flow Graphs:** P4Visor analyzes the pipelines of all P4 programs to be merged and identifies the tables to be merged using the algorithms presented in Section 4. Given this information, P4Visor merges the P4 programs by:

- (1) rewriting the Table IDs – to avoid conflicting IDs<sup>4</sup>,
- (2) rewriting “GoTo” statements for all tables except the merged tables to reflect the new Table IDs,

<sup>3</sup>Since we focus on merging different versions of the same program, and we anticipate a high degree of overlap between the parse graphs of the different programs.

<sup>4</sup>While each program may have unique Table IDs, multiple programs may reuse the same Table IDs which will create problems during compilations

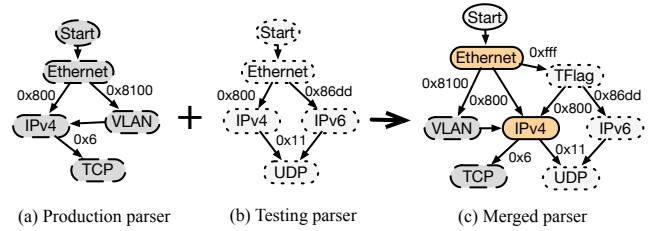


Figure 6: Merging two parse graphs.

- (3) for merged tables, P4Visor does one of two things: if the merged table leads to one table, then rewriting is obviously just rewriting the existing “GoTo” to use the appropriate ID; However, if a merged table leads to more than one table, *e.g.*, table “B” or “C” in Figure 3 (a), then P4Visor will add multiple “GoTo” statements, one for each branch.

Recall the above example, Figure 3, in the production control flow graph the next-hop for “C” is “D”, in the merged graph, table “C” will retain “D” as its default next-hop table modulo rewriting IDs to reflect D’s new ID; however, P4Visor will also add a “GoTo” that matches on the TFlag for testing control program and uses “A” as the next-hop for packets with the TFlag tag.

**Preserving Traffic Isolation:** Our framework must be able to offer isolation whenever it is required. Wherein, we define isolation as the following property: if two PDP Programs,  $P_1$ , and  $P_2$ , are isolated, then traffic for  $P_1$  is never processed by resources exclusively dedicated to  $P_2$ . Additionally, table entries controlled by one are never modified by the other. To enable isolation, we introduce an ACL-Bit (attached to tables) that provides access control overflow entries in the resource-sharing tables.

Observe that each table in the original P4 programs will map to exactly one table in the merged P4 program, while each table in the merged P4 program may correspond to one or more tables. We label all tables in the merged P4 Program that correspond to multiple tables in the original programs as *resource-sharing* tables. In Figure 3 (a), nodes C and B are both *resource-sharing* tables. For the *resource-sharing* tables, the P4Visor compiler will add the ACL-Bit to the table entries to provide traffic isolation for P4 programs. Combined with the TFlag which identify the packet as test packet, this ACL-Bit provides traffic isolation by allowing packets to match entries in the shared table only when packets match both the TFlag and the ACL-Bit. While the TFlag and ACL-Bit ensure isolation within resource-sharing tables, the TFlag alone ensures isolation between non-resource-sharing tables.

In addition to isolating the packet processing, the ACL-Bit also enables P4Visor to separate control over entries in the flow tables: The control plane for a P4 program can only modify the entries with the appropriate ACL-Bit value. Given this ACL-Bit, each P4 program can update the shared table correctly without side effects to the other P4 programs.

## 4 MERGING P4 PROGRAMS

In this section, we provide an overview of PDP-based resource constraints (Section 4.1), present the design of P4Visor’s source code merge model (Section 4.2), and conclude by theoretically analyzing the complexity and hardness of the problem (Section 4.3).

## 4.1 Background on P4 Compiler Constraints

In general, there are two kinds of constraints on a P4 program. These constraints are either placed on the compiler by the language (hardware target independent) or placed on the compiler by the hardware architecture (target-dependent). An example of a target-independent constraint is that there can be no loops in the control flow graph; hence, it needs to be a DAG. This constraint is invariant across all targets. However, the target-dependent constraints vary dramatically from target to target and are especially hard to enforce without intimate knowledge of the target hardware’s proprietary details. For example, RMT [7] has 32 stages in its pipeline while Intel’s FlexPipe [41] has 5 stages with different memory constraints for each stage.

To tackle these two constraints, P4 compilers are split into two components: a target-independent compiler (front-end compiler) and a target-dependent compiler (back-end compiler).

Prior work [27] has identified table size, program control flow, and hardware memory restrictions as the key issues faced by the P4 compiler.

In this paper, we focus on the design of target-independent merge optimizations. We aim to, first, provide a general optimization that benefits all hardware-targets. Our target independent optimization builds on the insight that merging different tables results in significant savings across all hardware targets for multiple reasons: merging tables reduces overheads associated with instantiating tables and merging tables results in large tables which take advantage of various hardware optimizations (we elaborate on this in Section 7). As part of future work, we will explore target-dependent optimizations.

## 4.2 Merging Optimization

Merging two P4 Programs is fundamentally equivalent to merging two weighted DAGs into a single weighted DAG with the added objective of minimizing space (*i.e.*, the # of nodes). To the best of our knowledge, no existing work has explored this problem: specifically, merging two weighted DAGs into one while maximizing overlap. The most closely related works [8, 43] provide suboptimal results, we elaborate on them in Section 9. Next, we more formally describe the problem.

We model a program’s control flow using a Table Dependency Graph (TDG) [27]  $G = (T, E)$  where vertices  $T = \{t_1, t_2, \dots, t_n\}$ <sup>5</sup> and edges  $E = \{(t_i, t_j) \mid t_i, t_j \in T\}$  map to the tables and the table dependency, respectively. Each table  $t_i \in T$  has three attributes:

- (1) the program id,  $t_i.pid$ , reflects the P4 program in which the table resides;
- (2) table ID,  $t_i.tid$ , reflects the table’s ID and helps to differentiate tables; and,
- (3) table size,  $t_i.size$ , reflects the memory footprint of the table (size is a function of width and number of entries defined).

Given  $G$ , we can compute the dependency matrix,  $D$ , of the graph as:  $D[t_i, t_j] = 1$  if there is a dependency path from  $t_i$  to  $t_j$  and  $D[t_i, t_j] = 0$  otherwise<sup>6</sup>.

<sup>5</sup>  $n$  is the total number of tables in the pipeline.

<sup>6</sup>  $D[t_i, t_i] = 0$  because P4 programs are generally acyclic graphs.

For simplicity, we formalize the PDP-merge problem for two P4 programs, but the problem formulation and analysis generalizes to cases with more than two P4 programs.

**Objective:** Our goal is to merge two programs – a production version denoted as  $G_r = (T_r, E_r)$  with the dependency matrix  $D_r$  and a testing version  $G_s = (T_s, E_s)$  with the dependency matrix  $D_s$  into a single program  $G_m = (T_m, E_m)$  with the dependency matrix  $D_m$ , while minimizing the total resources required. In this paper, we only focus on table memory resources. Restated, our object is maximizing sharing resources:

$$\max \sum_{i=1}^{|T_m|} w_i \quad (1)$$

where  $w_i$  is the weighted contribution of reducing the resources in  $G_m$  used by table  $t_i \in T_m$ .

We define the set of *resource-sharing* tables,  $T_{ms}$ , as a subset of tables in the merged TDG  $G_m$ : these tables in  $T_{ms}$  are merged from multiple tables in the original programs, which satisfy the following constraints: equivalence, correctness, and loop-freedom.

For each table  $v_i \in T_{ms}$ ,  $w_i$  captures both the memory type and table size. Currently, the memory size is calculated as a function of the number of entries and the width of each entry:

$$w_i = c_i \cdot len_i \cdot width_i$$

where  $len_i$  and  $width_i$  are the number of entries and width of an entry in table  $t_i$  respectively.  $c_i$  is a configurable weighted coefficient that allows an administrator to guide our optimization algorithm to merge tables that the administrator cares about. For example, if an administrator only cares about the TCAM tables, she can set the table weights of all non-TCAM types to 0. As a preprocessing step, P4Visor sets  $w_i = 0$  for each table  $v_i \notin T_{ms}$  which shares no table resources with other tables because these tables cannot be merged.

Note that when the weights for all tables  $T_{ms}$  are equal, the objective function (1) leads to a merged TDG that minimizes the total number of tables.

**Target-Independent Constraints:** Two tables,  $t_{r_i} \in T_r$  and  $t_{s_j} \in T_s$ , can be merged if and only if three constraints are satisfied:

- (i) *Equivalence:* The two tables are structurally equivalent (same actions and match fields but they can vary in the number of declared entries). Here the equivalent tables are assigned the same id, that is  $t_{r_i}.tid = t_{s_j}.tid$ .
- (ii) *Correctness:* the table dependencies of both tables are maintained – correctness is preserved.

$$\begin{cases} D_m[t_{r_i}, t_{r_j}] = D_r[t_{r_i}, t_{r_j}], \forall t_{r_i}, t_{r_j} \in T_r \\ D_m[t_{s_i}, t_{s_j}] = D_s[t_{s_i}, t_{s_j}], \forall t_{s_i}, t_{s_j} \in T_s \end{cases} \quad (2)$$

- (iii) *Loop-free:* the resulting graph is loop free, that is, the dependency matrix of  $G_m$  satisfies  $\forall t_i, t_j \in T_m$ ,

$$D_m[t_i, t_j] \cdot D_m[t_j, t_i] = 0 \quad (3)$$

**Target-Dependent Constraints:** While this work focuses on target-independent constraints, here, we briefly sketch out how target-dependent constraints can be introduced into our problem formulation.

Abstractly, we can introduce target-dependent constraints by introducing hardware information. One constraint placed by hardware is the number of physical stages. For example, RMT [7] has 32 stages, and thus RMT can only support P4 programs whose crucial dependency path length is no more than 32. To overcome this limitation, we can add a constraint that limits the merged TDG's critical path length to less than 32. This may force our algorithm to explore solutions that create merged programs that do not maximize overlaps, but that ensure shorter critical dependency paths.

As part of future work, we will study the constraints of dominant PDP hardware targets and incorporate them into our algorithm.

### 4.3 Complexity Analysis

Our TDG merging problem can be reduced to and from the Maximum Weighted Independent Set (MWIS) problem: a problem which has been proven to be NP-Complete [18]. In this section, we provide a sketch of how to reduce our problem to and from the MWIS problem.

We define a function  $v(m, i)$  which returns the table from TDG  $G_m$  whose table ID is  $i$ , thus,  $v(m, i) = t_{m_i}$  and  $t_{m_i}.tid = i$ . To do this reduction, we define a merge candidate set,  $T_p$ , as the set of all tables in  $G_r$  and  $G_s$  that satisfy the equivalence requirement defined in constraint (i).

By definition of constraints (i) to (iii) in Section 4.2, all the tables in production and testing programs follow the Lemma 4.1 (proved in the appendix).

LEMMA 4.1.  $\forall t_i, t_j \in T_{ms}$ ,  $v(s, j)$  and  $v(s, i)$  have

$$D_r[v(r, i), v(r, j)] \cdot D_s[v(s, j), v(s, i)] = 0 \quad (4)$$

Next, let us construct a new undirected graph  $G_p = (T_p, E_p)$  where the vertex set of the graph is  $T_p$  and the edge set of the graph is  $E_p$ . Given this definition, we define  $\forall t_i, t_j \in T_p$ ,

$$E_p[t_i, t_j] = \begin{cases} 1 & D_r[v(r, i), v(r, j)] \cdot D_s[v(s, j), v(s, i)] = 1 \\ & \text{or } D_r[v(r, j), v(r, i)] \cdot D_s[v(s, i), v(s, j)] = 1 \\ 0 & \text{Otherwise} \end{cases} \quad (5)$$

Taken together, formulas (4) and (5) provide us with a way to formally reason about the relationship between  $T_{ms}$  and  $G_p$ . Lemma 4.2 (proved in the appendix) provides this relationship.

LEMMA 4.2. *The set of resource-sharing tables  $T_{ms}$  is a subset of vertices in graph  $G_p$ , no two of which are adjacent, that is,  $\forall t_i, t_j \in T_{ms}$ ,*

$$E_p[t_i, t_j] = 0 \quad (6)$$

**Reducing PDP-Merge to MWIS:** Lemma 4.2 restated shows that analyzing graph  $G_p$  to identify the set of tables  $T_{ms}$  can be reduced to the independent set problem in polynomial time of  $O(|T_p|^2)$ . Essentially, in constructing  $G_p$ , we only keep the dependencies in both  $D_r$  and  $D_s$  that provide the forward and reverse direction between two nodes. Take nodes  $A, D$  in Figure 3 as an example, there is a dependency from node  $A$  to  $D$  in one program as well as a dependency path from  $D$  to  $A$  in another program. We keep these types of forward and reverse dependencies when creating  $G_p$  and delete all others dependencies.

To satisfy our objective of maximizing the shared table resources, we need to find the maximum weighted independent set in graph  $G_p$ , known as MWIS problem, an NP-Complete problem [18].

**Reducing MWIS to PDP-Merge:** Next, we show how to reduce a given MWIS problem to our merging problem. The key lies in transforming a given weighted undirected graph  $G_p = (T_p, E_p)$  in the MWIS problem to two weighted DAGs,  $G_r$  and  $G_s$ , to be merged with the objective of maximizing the weights of the final DAG. More specifically, we can construct the dependencies matrix of two DAGs from  $G_p$  as follows:

$$D_r[i, j] = \begin{cases} E_p[i, j] & \text{if } i > j \\ 0 & \text{Otherwise} \end{cases} \quad (7)$$

$$D_s[i, j] = \begin{cases} E_p[i, j] & \text{if } i < j \\ 0 & \text{Otherwise} \end{cases} \quad (8)$$

where  $i, j = 0, 1, 2, \dots, |T_p|$  are the indices of the nodes in graph  $G_p$ . We set  $T_{ms}$  as a feasible independent set of  $G_p$ . Similarly, with lemma 4.2, we know that  $T_{ms}$  is a feasible set of *resource-sharing* tables when merging two constructed DAGs  $G_r$  and  $G_s$ . Further, as each node has a weight, solving the maximum weighted independent set of  $G_p$  is equal to the identification of the set of tables with maximum shared resource when merging  $G_r$  and  $G_s$ .

Thus, the maximum weighted independent set (MWIS) problem, an NP-Complete problem [18], can be reduced to our problem in polynomial time  $O(|T_p|^2)$ . That is to say, merging two weighted DAGs into one weighted DAG with the objective of maximizing weights is an NP-Complete problem.

## 5 EFFICIENCY

Next, we design a heuristic to efficiently solve the problem in real time (Section 5.1) and discuss a systematic approach for configuring resource sharing of entries within the merged tables (Section 5.2).

### 5.1 P4Visor Heuristic Merging

A naive approach for solving the “merge” problem is to perform a brute-force search through all potential combinations in  $G_m$  to find the solution which provides the maximum overlap: the best-known optimal algorithm for solving the maximum independent set problem is Bron-Kerbosch. We implemented the extended Bron-Kerbosch [24] and observed that it can only handle small graphs and was unable to scale to large graphs (i.e., greater than 80 nodes): In particular, given a 7-day time limit, we were unable to solve the Bron-Kerbosch algorithm for graphs with over 80 nodes. Thus, Bron-Kerbosch was unable to process the largest DAG in our dataset (Switch.P4 which has over 120 tables). Motivated by the inadequacies with Bron-Kerbosch, we designed a new heuristic to solve the merge problem.

**Heuristic** Our heuristic is based on simulated annealing (SA) which has proven effective in solving the MWIS problem [1].

In our heuristic, each state of the search space is defined as a subset  $V_{sub}$  of the vertex set of graph  $V_p$  and every vertex in  $V_{sub}$  is nonadjacent to the other vertices. Motivated by prior work [4, 40], our heuristic generates neighboring states to explore using one of the following two procedures:

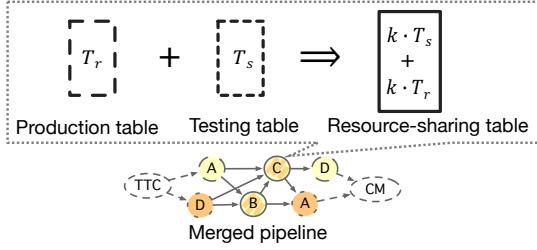


Figure 7: Controlling Resource Sharing

- (1) adding one vertex  $v_i \in V_p \setminus V_{sub}$  to  $V_{sub}$  and deleting all the vertices,  $v_j$  in  $V_{sub}$ , that are adjacent to  $v_i$ .
- (2) adding two nonadjacent vertices  $v_i, v_j \in V_p \setminus V_{sub}$  to  $V_{sub}$  and deleting one vertex from  $V_{sub}$  that is adjacent to both  $v_i$  and  $v_j$ .

The energy function, or evaluation function, of our heuristic is defined as the total weight,  $E(V) = \sum_{v_i \in V} weight(v_i)$ , of the vertices in current search state,  $V$ , where  $weight(v_i)$  is the weight of each vertex  $v_i \in V_{sub}$ . A new state will be accepted if its energy, *i.e.*,  $E(V_{n+1})$ , is larger than the current state's energy  $E(V_n)$ ; otherwise we accept the new state with probability of  $e^{-\frac{\Delta E}{t}}$ . The temperature,  $t$  is initially set to 100 and decreases linearly to 1 at each iteration based on this equation  $t(n+1) = t(n) * 0.99$ . We terminate the search when the temperature decreases to 1.

A well known problem of simulated annealing-based heuristics is that they can often get stuck in local optimas. To avoid this problem, we run the simulated annealing process many times which increases coverage over the search space and introduces more randomness. The increased randomness and larger coverage over the search space, enables our heuristic to avoid local optimas.

## 5.2 Controlling Resource Sharing

In general, merging the control flow graphs of two P4 programs consists of two major steps. The first is to identify the tables to be merged in both P4 programs,  $G_r$  and  $G_s$ , that satisfy the ‘‘correctness’’ constraints. The second step is to merge the control flow according to the identified tables.

In merging two tables,  $t_{r_i}$  and  $t_{s_j}$ , the resulting table,  $t_{m_x}$ , can vary in size, # of entries, ranging from  $\max(t_{r_i}.size, t_{s_j}.size)$  (the merged tables reuse 100% of their resources) to  $t_{r_i}.size + t_{s_j}.size$  (the merged tables do not reuse any resources and the logical size of the merged table is equivalent to the sum of the original tables).

At both extremes, P4Visor provides benefits. At the extreme where no table entries are shared, *i.e.*,  $t_{m_x}.size = t_{r_i}.size + t_{s_j}.size$ , the merge provides benefits because it enables  $G_m$  to fit within smaller memory by reducing the overheads associated with instantiating individual tables in hardware, *e.g.*, in Xilinx's Virtex-7 FPGAs [55] instantiating a TCAM table with 64 bits width  $\times$  256 bits depth consumes 2 RAM blocks but instantiating TCAM tables with 4 times (128 bits width  $\times$  512 bits depth) or 16 times (256 bits width  $\times$  1024 bits depth) more memory consumes only 3 or 5 RAM blocks respectively. Thus, merging would save resources, even if the merged table contains the same number of entries as the original tables combined.

At the other extreme, where switch memory is reused, that is,  $t_{m_x}.size = \max(t_{r_i}.size, t_{s_j}.size)$ , then precious switch memory is being saved by sharing resources across tables. This savings is in addition to the savings of overheads described earlier. At this extreme, the table resources are shared proportionally between the different programs based on the fraction of traffic allocated to each program.

To explore this trade-off, P4Visor exposes a parameter,  $k$ , to the operator through the P4Visor Interface. This parameter allows the operator to control the amount of sharing:  $k = 1$  means no sharing while a  $k = 0$  means proportional sharing. Figure 7 demonstrates the controlling resource sharing.

## 6 P4VISOR IMPLEMENTATION

We have implemented the P4Visor compiler in 3000+ lines of Python code and 800+ lines of C++ code. The controller application, P4Visor interface and P4Visor management agent are all developed with Python in over 100+ LoCs. The P4Visor compiler takes as input the high-level intermediate representations of P4 programs (*i.e.*, HLRIR) and merges them into one program. Merging the high-level IR allows us to operate at a platform independent level while maintaining the complete semantics of the P4 language. Currently, we only support merging of P4<sub>14</sub> programs. As part of future work, we will extend P4Visor to support P4<sub>16</sub> programs. P4Visor's source code is online in our Github repository [6].

### 6.1 Supporting Flexible Testing Operators

We now discuss, the implementation details of several interesting components: specifically, the Differential Testing specific module. In this discussion, we also, demonstrate the flexibility of the testing operators provided by P4Visor.

**Testing Traffic Control (TTC):** In Figure 8, we present a code excerpt from our implementation of the TTC component. Recall, the key goal of the TTC is to provide flexible control over the sampling and selection of traffic for testing. P4Visor provides both runtime and compile-time control which allows the administrator to alter sampling configuration program: at compile time the administrator can configure more static aspects of the flow spec to match on and at runtime, the administrator can specify the exact values to sample on.

The control flow of TTC is implemented in table `testing_traffic_control` (Lines #11-14), the TTC uses the `set_testingbit` variable to determine if a packet should use the testing pipeline or the production pipeline.

**Run Time Configuration:** The administrator can control the sampling rate (by changing the registers (Line #1)) and the subnets to be sampled (by modifying the entries in `testing_traffic_identify` table (Line #8)). The sampling frequency is implemented as a special action, `sample_testing_pkt` (in Lines #2-5), which uses two registers, `cnt` and `Rate`, to determine the sampling rate, *e.g.*, sample one packet in every  $R$  packets.<sup>7</sup> The subnet sampling is performed by comparing the addresses in the packet against the entries in table `testing_traffic_identify` (Line #8).

<sup>7</sup>Due to limited arithmetic operations supported by P4, the TTC can only support the sample ratio of  $1/R$  ( $R = 1, 2, 3, \dots$ ).



*Compile Time Configuration:* Additionally, the structure of table `testing_traffic_identify` can be altered, through the PVI, to re-configure the TTC and allow the TTC to match packets for sampling based on other aspects of the FlowSpec beyond subnet.

```

1 registers cnt, Rate
2 action sample_testing_pkt() {
3   register_write(cnt, 0, cnt+1);
4   modify_field(testing_meta.testingbit, cnt%Rate);
5 }
6 table testing_traffic_identify {
7   //fields are configurable using P4Visor Interface
8   reads {ipv4.dstAddr : lpm;}
9   actions {sample_testing_pkt; set_testingbit;}
10 }
11 table testing_traffic_control {
12   reads {testing_meta.testingbit : exact;}
13   actions {goto_test_pipe; goto_prod_pipe}
14 }

```

Figure 8: Code Excerpt from our TTC Implementation.

**Comparator:** In Figure 9, we present an pseudocode for the Comparator Module. The Comparator is implemented using a set of flow tables with compound actions. The output of each version of the program are recorded in a set of metadata (i.e., *meta\_p*) and then compared by the Comparator (Line #2) to determine if the versions are different. If a difference is detected, the action `diff_procedure` is used to create a packet to send to the controller. To overcome a limitation of our target platforms, we create a new packet by multicasting the original packet and sending a version to the controller (Line #6-9).

```

1 //compare the outputs
2 if(testing_meta.meta_p!=testing_meta.meta_t){
3   apply(diff_procedure);
4 }
5 //an example procedure configuration
6 action diff_procedure(testing_meta, mcast_group) {
7   update_fields(testing_tag, testing_meta);
8   set_output_mcg(mcast_group);
9 }

```

Figure 9: Pseudocode for Comparator.

## 6.2 Limitation of Existing PDP Targets

PDPs are expected to provide a rich set of packet processing features, e.g., the action primitives defined in P4. However, current PDP targets, e.g., software switch Bmv2 [13] or FPGA-based hardware from Xilinx SDNet [54], can only support a limited set of P4’s features. Several of the key P4 features required to enable P4Visor include (1) stateful registers; (2) packet cloning, for creating multiple copies of a packet to be processed by different programs; and (3) in-switch packet generator, for generating and sending a packet to the controller that summarizes differences between the two P4 programs.

- While packet cloning primitives are defined in the P4 specification, the clone feature is not supported by the Bmv2 target. We address this problem by attaching attributes of the packet to the metadata and recirculating the packet and metadata through the pipeline for processing by the alternative programs. Thus, by recirculating the packet, multiple versions of a PDP Program can independently process the same packet.

- In-switch packet generator is not supported by either the Bmv2 or FPGA targets (Xilinx SDNet). To send the outcome of testing to the controller, P4Visor adds those fields to the pre-configured TFlag, inserts the TFlag to a copy of the packet, and then sends a copy of the packet out to the controller.<sup>8</sup>
- Our hardware target is even less flexible than the Bmv2 due to the limitation of the current development toolchain (SDNet). Specifically, SDNet does not support stateful register which impacts our design of the comparator and limits the set of programs we can deploy. To support P4Visor on our FPGA-based hardware target, we have implemented those primitives in low-level hardware (i.e., the Testing Traffic Control module is implemented with 1000+ lines of Verilog code). We believe these hardware limitations will be addressed with the evolution of the SDNet toolchain.

## 7 EVALUATION

### 7.1 Experiment Setup

We have evaluated P4Visor on both a software (Bmv2 [13]) and a hardware (ONetSwitch [23]) programmable data plane.

**Software PDP:** On the Bmv2 target, we analyzed the following programs: Reference Switch.P4 [14], L2 switch, Simple Router, NAT, VLAN and Arp-Proxy, Flowletting [15], and Heavy Hitters [48]. The Bmv2 runs in mininet with a single switch, two hosts for testing, and a third host for running the controller. Before testing, we install flow entries into the tables so that the two end hosts can ping each other.

**Hardware PDP:** On the ONetSwitch target, we were only able to evaluate the following programs: L2 Switch, Simple Router, and VLAN. We were limited in the set of programs evaluated because ONetSwitch builds on Xilinx’s Zynq SoC [23] which only supports a subset of P4’s language features (see Section 6.2). To test the performance of the switch, we connect two PCs with 10G NIC to the ONetSwitch45 switch, due to NIC limitations, the maximum achievable throughput for our servers is 5Gbps. We used iPerf to generate traffic between the hosts and similarly crafted rules to force traffic through as many tables as possible.

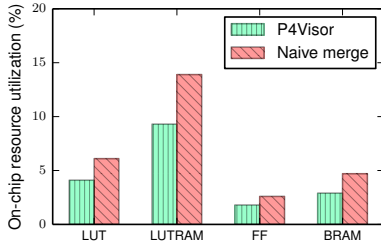
### 7.2 Performance Benefits and Overheads

Here, we evaluate the overheads of P4Visor and analyze the practical benefits of source code merging as a lightweight virtualization primitive.

*7.2.1 Benefits of Resource Sharing.* To understand and quantify the benefits of resource sharing, we have compared P4Visor’s merge algorithm against a *Naive merge* algorithm [44], which is a greedy algorithm for MWIS problem. When solving the problem, *Naive merge* selects a vertex of minimum degree, removes it and its neighbors from the graph until no vertex available.

Our results show that merging introduces significant benefits for three distinct reasons: First, instantiating a table into hardware incurs some overheads. Thus by having two programs sharing a table, we ameliorate the associated overheads and this translates

<sup>8</sup>Recall, the TFlag is removed at the edge switch and thus the endhosts never receive the TFlag.



**Figure 10: Impact of program merging approaches on memory utilization.**

into memory savings. For example, while multiple tables may use the same actions, these actions need to be independently stored for each table and by merging tables, we reduce the number of instances of these actions. Second, as tables grow in size, several of the resources increase in a sub-linear fashion due to hardware-specific optimizations, *e.g.*, the BRAM memory in Xilinx includes optimizations that result in sub-linear growth. Third, when we modify the parameter  $k$  which impacts the amount of sharing, we reduce the total number of entries. This introduces yet more savings.

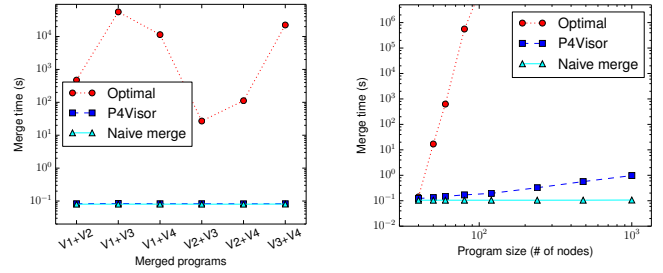
To illustrate the first two points, we analyze a simple P4 program (the router program [5]) consisting of two tables: a TCAM IPv4 routing table (32 bits width, 256 entries) and an exact match IPv6 routing table (128 bits width, 256 entries); each table has two actions. For this analysis, we set  $k = 1$  which means both the P4Visor-merged and the naive-merged programs will have the same number of entries. We merge two router programs into one and then compile the merged programs to the ONetSwitch45. In our analysis, we focus on the four different kinds of memory resources of the Xilinx chipset: LUT, LUTRAM, FF, and BRAM.<sup>9</sup> Figure 10 illustrates our first point – our heuristic merging with P4Visor results in a 32% to 49% savings in resources compared with the naive merging algorithm. To illustrate the second point, we analyze the amount of resources required to support tables of varying sizes. We observe that while most resources grow linearly with the size of the tables, the BRAM grows sub-linearly (figure omitted due to space limitations).

**Takeaway.** Intelligently merging tables leads to tremendous resource savings. We anticipate these savings to only grow as P4 programs become even more complex.

**7.2.2 Performance Overheads.** To evaluate the performance overheads of P4Visor, we randomly select two of the evaluated PDP programs, *i.e.*, L2switch, Router, VLAN, NAT, and Flowleting (only the first three for hardware switch), merge them with P4Visor, and compare the throughput/latency of running traffic through the merged program against that of running traffic through the unmodified programs – we compare the merged program against the better performing of the two original PDP programs.

Our experiments, not shown due to space constraints, demonstrate that P4Visor introduces minimal overheads. Specifically, the TTC and the Comparator modules which add tags to packets and perform comparisons introduce minimal overheads. In the software switch, the throughput decreases by less than 1.5% and the delay

<sup>9</sup>The LUT, LUTRAM, and BRAM are mainly used to store Table structures – with BRAM used to store large TCAM tables. The FF, is however, mainly used to store timing and control signals.



(a) Real P4 Programs

(b) Synthetic P4 Programs.

**Figure 11: Runtime of program merging approaches.**

penalty is less than 3%. For the hardware switch, the overhead is much smaller, both the throughput and delay overheads are less than 1%.

**Takeaway.** P4Visor introduces several tables and actions to support the different testing paradigms; however, these constructs introduce minimal performance overheads to the network (less than 1% in hardware) making them highly desirable for today’s networks.

### 7.3 Analytical Evaluation of the Heuristic

To evaluate the efficiency and accuracy of our heuristic, we compare our heuristic with the optimal solution (Bron-Kerbosch) and a naive greedy merge [44]. Note: we were unable to evaluate the optimal approach on programs with over 80 nodes because the optimal algorithm failed to provide a solution. In these evaluations, we focus on two kinds of P4 programs:

First, real P4 programs, is based on the reference Switch.P4 [14], which contain 82 tables in the ingress pipeline and 41 nodes in the egress pipeline. As Switch.P4 is built in a configurable fashion, we create different versions by turning on or off specific functionality. Specifically, we created the following four versions:

- Switch.P4-V1: by turning on the OpenFlow processing module. It has 84 ingress nodes and 53 egress nodes.
- Switch.P4-V2: by turning off the Tunnel processing module. It has 66 ingress nodes and 30 egress nodes.
- Switch.P4-V3: by turning off the ACL processing module (MAC, IPv4, IPv6, RAACL/PBR). It has 76 ingress nodes and 37 egress nodes.
- Switch.P4-V4: by turning off the Multicast processing module. It has 73 ingress nodes and 39 egress nodes.

Second, synthetic P4 programs, which enable us to systematically evaluate the accuracy and efficiency of our heuristic at scale – larger than the largest known P4 program (Switch.P4). We generate synthetic programs ranging in size from 30 to 1000 tables with randomly generated dependencies (edges) – we generate the edges to ensure that the graph maintains a graph density of 0.4.

**Efficiency** To evaluate the efficiency of our heuristic, we measure the time it takes to merge two randomly selected P4 programs. Figure 11(a) presents the runtimes for merging real P4 programs. We observe that our heuristic and the naive algorithm consistently take a similar amount of time (0.1 seconds) and both are considerably faster than the optimal algorithm. Figure 11(b) presents the

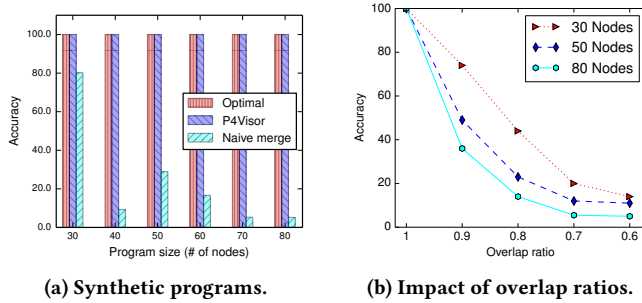


Figure 12: Accuracy of program merging approaches.

runtimes for merging synthetic P4 programs: this figure highlights the relationship between the runtime of the three approaches and P4 programs size (in # of nodes): the optimal algorithm shows an exponential growth, while both our heuristic and the greedy approach show a linear growth as a function of program size. With the larger graphs, we observe that while our heuristic performs slower than greedy, it’s performance is still acceptable.

**Accuracy** Next, we analyze the accuracy of the different approaches to understand the cost of the performance improvements. To evaluate the accuracy, we compare the solutions generated by the different approaches against the optimal solution. Within the real programs (not shown due to space constraints), we observed that both the greedy approach and our heuristic achieved 100% accuracy; however, for the synthetic programs (Shown in Figure 12(a)) we observed that the greedy approach achieved an average accuracy of 30% while our heuristic was able to achieve 100% in all situations.

Upon further analysis, we observed that the accuracy of the greedy heuristic is a function of the ratio of overlap between the different programs. Recall, the real programs are all variants of Switch.P4 thus we expect there to be significant overlaps. If the overlap is extremely high, as it is common when minor changes are made, then the greedy heuristic performs well; however if the overlap is low, e.g., when significant changes are made, then accuracy drops. To illustrate this point, in Figure 12(b) we explore the impact of overlap ratio on accuracy. From this figure, we observe that for the greedy approach, the amount of overlap has a large impact on its accuracy.

**Takeaway.** While our heuristic is slower than the greedy approach, our heuristic scales linearly and provides better accuracy across a broader set of scenarios. In short, our heuristic is fast and accurate.

## 7.4 Use Case: Testing P4 Programs

In section 6.1, we showed how to flexibly configure the fields to be tested and actions to be performed using P4Visor interface. In this section, we demonstrate the use of P4Visor to perform testing and illustrate how these interfaces may be configured. Specifically, we use P4Visor to perform Differential Testing to test the behaviors of two versions of the P4 Router program. Unlike the previous sections which focus on overheads and accuracy, here we explore the operational interactions involved with using P4Visor.

**Testing Setup:** To use P4Visor, we (1) configured P4Visor to record and compare the 32-bits next-hop metadata fields of the programs. To handle the detected differences, (2) configured the Comparator to send the packets along with the outputs from two programs to the controller (the same as the configuration in Figure 9), and (3) fed the configuration files and the two P4 programs into P4Visor’s interface to produce the merged program.

**Testing results:** We evaluated the merged program on the Bmv2 target. At runtime, we control the routing tables of the two programs with two different routing applications. We observed that the P4Visor application can detect differences, via the P4Visor-generated messages, within milliseconds once the control planes install different routing entries for the same flow. With the help of the outputs stored in the messages, an administrator can further debug and analyze the behavior of the tested programs.

## 8 RELATED WORKS

The most closely related work [30] explore source code merging as a method for providing virtualization. We explore a similar approach but in the P4 domain and tackle a host of domain-specific issues. Moreover, we prove the complexity of the merge problem. Below we explore related works on SDN composition, DAG-Merging, and other recent work on programmable data planes.

**PDP:** Many have explored virtualization [2, 26, 45], update [29, 35, 37], and state-management [50, 58] techniques for traditional SDNs. Despite the growing emergence of PDP-based architectures and solutions, to-date, there are few principled approaches for supporting testing P4 Programs.

Most work focus on applications of PDPs [11, 16, 32, 36] or developing interfaces and primitives to enrich existing PDPs environment [7, 34, 47]. Our work falls in the latter class and argues for a principled extension of PDPs to include interfaces, abstractions, and primitives to enable testing — in short, to support rapid prototyping.

**PDP Compiler:** Several works [19, 27, 53, 54] have explored challenges associated with compiling P4 programs to various hardware targets, e.g. FPGA [53, 54] or RMT [19, 27]. Our work can benefit from these approaches by using these approaches within the back-end compiler – note: in this paper, we focus our emphasis on front-end optimizations.

**PDP Virtualization:** While related work [22, 56, 57] proposed a general virtualization primitive for P4, P4Visor provides lighter weight virtualization primitive for testing based on code merging.

**SDN Composition:** Orthogonal work on composition and modularity in SDNs [9, 17, 21, 26, 38] focus on SDN rules and not on P4 program’s source code. Concurrent work on composition [49] aims to support orthogonal composition operators. In this work, we present the first attempt to formalize the problem and present a framework with supporting algorithms and abstractions to enable composition within PDPs.

**DAG-Merging:** The problem of merging two DAGs has been explored by several others [8, 28, 43]. At a high level, our work differs in two ways: first, we investigate a different formulation

objective and second, we prove the complexity of the DAG merge problem and provide a more efficient heuristic based on simulated annealing. Below, we elaborate on these differences.

OpenBox [8] and SNF [28] merge processing graphs for network functions (NFs) with the objective of minimizing the path length (reducing packet processing latency in the merged processing graph). Our goal differs from prior work because we aim to maximize the number of merged nodes (minimizing the resources used by the merged graph). A significant implication of this difference between our objectives is that while our heuristic always reduces the size of the merged graph, the algorithms presented by OpenBox and SNF may increase the number of nodes in an attempt to minimize path lengths.

Saha et al.[43] provide a sub-optimal heuristic for merging two unweighted DAGs into one unweighted DAG with the objective of minimizing the number of vertices in the final DAG. Our objective varies: we want to simultaneously decrease the number of vertices while maximizing the magnitude of overlap since the vertices in our problem are weighted. Saha et al. show that merging DAGs is the dual of finding the maximum length longest common subsequence between pairs of topological sorted DAGs. Unfortunately, while this problem can be solved in  $O(V^2)$  time, the efficiency of the solution is highly dependent on the topological sort of the DAG and can be highly ineffective. Instead, we present an efficient greedy heuristic based on simulated annealing.

## 9 DISCUSSION

**Control Plane Overheads.** In addition to the data plane overheads, P4Visor introduces overheads to the control plane. For example, the P4Visor agents running on the controller and the switch PDP devices, have to multiplex and demultiplex messages between the controller to the local P4 programs. This translation introduces processing overheads and also memory overheads because the agents need to maintain a mapping and perform the translations. Additionally, re-using the control channel between the controller and switches to transfer packets summarizing the result of the tests reduces the available bandwidth on the control channel. As part of future work, we plan to explore approaches, e.g., SwitchVisor [12], to effectively share and partition these control plane resources.

**Target Dependent Optimizations** As discussed in Section 4, our current efforts focus on target-independent optimizations (i.e., front-end compiler), as part of future work we will extend our formulation to tackle the back-end compiler by introducing constraints and objectives specific to the hardware targets.

**Seamless Reconfiguration** While full virtualization provides support for headless updates (reconfiguring the data plane without a reboot), our approach requires a reboot after every reconfiguration. As part of future work, we plan to tackle issues related to these reboots by intelligently migrating state, e.g., with SwingState [34], and reconfiguring paths, e.g., with zUpdates [33], during the reboot to eliminate disruption.

**Composition Operators** This work has focused on supporting testing-specific composition operators; however, as part of on-going work we are exploring composition operators for enabling code

modularity, e.g., parallel and sequential composition. Supporting these operators requires extending our current formulations to account for operator specific constraints.

## 10 CONCLUSION

In this paper, we propose a lightweight virtualization primitive for testing P4 programs through code merging. To support this primitive, we present a framework, called P4Visor, which uses compiler optimizations and program analysis to achieve efficient source code merging. We evaluate the theoretical complexity of the merging algorithm and present an efficient greedy heuristic. Our work opens up space for implementing a wide range of composition operators and frameworks for P4 programs.

## ACKNOWLEDGMENTS

We thank our shepherd Eric Keller, and the anonymous CoNEXT reviewers for their invaluable comments. This work is supported in part by the National Key Research and Development Program of China (2017YFB0801703), the National Science Foundation (through grants CNS-1749785 and CNS-1819109), and the NSFC (61672425, 61702407).

## A APPENDIX

### A.1 Proof of Lemma 4.1

We can proof lemma 4.1 by contradiction. Let us assume that  $\exists t_i, t_j \in T_m$  so that

$$D_r[v(t, i), v(t, j)] \cdot D_s[v(s, j), v(s, i)] = 1$$

then we know  $D_r[v(t, i), v(t, j)] = 1$  and  $D_s[v(s, j), v(s, i)] = 1$ . As  $i, j$  are the ids of the merged tables satisfying the table dependency consistency, according to *Rule1* we have

$$D_m[v(m, i), v(m, j)] = D_r[v(t, i), v(t, j)] = 1$$

$$D_m[v(m, j), v(m, i)] = D_s[v(s, j), v(s, i)] = 1$$

which means the merging of tables  $v(t, i), v(s, i)$  and the merging of  $v(t, j), v(s, j)$  introduce a dependency loop to the merged graph  $D_m$ . By *Rule2*,  $D_m$  is loop free. This is a contradiction. QED.

### A.2 Proof of Lemma 4.2

We can proof lemma 4.2 by contradiction similar with the proof of lemma 4.1. Assume that  $\exists t_i, t_j \in T_m$  so that  $E_p[t_i, t_j] = 1$ , then according to equation (5) we can get  $D_r[v(t, i), v(t, j)] = 1$  and  $D_s[v(s, j), v(s, i)] = 1$ . This will lead to the same contradiction shown in the proof of lemma 4.1. Hence, we have  $\forall t_i, t_j \in T_m, E_p[t_i, t_j] = 0$ . QED.

## REFERENCES

- [1] Emile Aarts and Jan Korst. 1989. *Simulated Annealing and Boltzmann Machines: A Stochastic Approach to Combinatorial Optimization and Neural Computing*. John Wiley & Sons, Inc., New York, NY, USA.
- [2] Ali Al-Shabibi, Marc De Leenheer, Matteo Gerola, Ayaka Koshibe, Guru Parulkar, Elio Salvadori, and Bill Snow. 2014. OpenVirteX: Make Your Virtual SDNs Programmable. In *Proceedings of the Third Workshop on Hot Topics in Software Defined Networking (HotSDN '14)*. ACM, New York, NY, USA, 25–30. <https://doi.org/10.1145/2620728.2620741>
- [3] Richard Alimi, Ye Wang, and Y. Richard Yang. 2008. Shadow Configuration As a Network Management Primitive. In *Proceedings of the ACM SIGCOMM 2008 Conference on Data Communication (SIGCOMM '08)*. ACM, New York, NY, USA, 111–122. <https://doi.org/10.1145/1402958.1402972>



- [4] Diogo V. Andrade, Mauricio G. C. Resende, and Renato F. Werneck. 2012. Fast local search for the maximum independent set problem. *Journal of Heuristics* 18, 4 (01 Aug 2012), 525–547. <https://doi.org/10.1007/s10732-012-9196-4>
- [5] The Authors. 2018. The P4 Router Programs. <https://github.com/Brown-NSG/P4Visor/tree/master/FPGAtarget/p4program>. (2018).
- [6] The Authors. 2018. The P4Visor Compiler for BMV2 target. <https://github.com/Brown-NSG/P4Visor>. (2018).
- [7] Pat Bosshart, Glen Gibb, Hun-Seok Kim, George Varghese, Nick McKeown, Martin Izzard, Fernando Mujica, and Mark Horowitz. 2013. Forwarding Metamorphosis: Fast Programmable Match-action Processing in Hardware for SDN. *SIGCOMM Comput. Commun. Rev.* 43, 4 (Aug. 2013), 99–110. <https://doi.org/10.1145/2534169.2486011>
- [8] Anat Bremner-Barr, Yotam Harchol, and David Hay. 2016. OpenBox: A Software-Defined Framework for Developing, Deploying, and Managing Network Functions. In *Proceedings of the 2016 ACM SIGCOMM Conference (SIGCOMM '16)*. ACM, New York, NY, USA, 511–524. <https://doi.org/10.1145/2934872.2934875>
- [9] Marco Canini, Daniele De Cicco, Petr Kuznetsov, Dan Levin, Stefan Schmid, Stefano Vissicchio, et al. 2014. STN: A robust and distributed SDN control plane. *Open Networking Summit* 490 (2014).
- [10] Inc. Cavium. 2018. XPliant Ethernet Switch Product Family. (2018). <http://www.cavium.com/XPliant-Ethernet-Switch-ProductFamily.html>
- [11] Ed Doe Changhoon Kim, Parag Bhide. 2016. In-band Network Telemetry (INT). <http://p4.org/wp-content/uploads/INT/INT-current-spec.pdf>. (2016).
- [12] Huan Chen and Theophilus Benson. 2017. Switch-visor: Towards Infrastructure-level Virtualization of SDN Switches. In *Proceedings of the 2Nd Workshop on Cloud-Assisted Networking (CAN '17)*. ACM, New York, NY, USA, 25–30. <https://doi.org/10.1145/3155921.3158431>
- [13] P4 Language Consortium. 2017. P4 software switch (behavioral model) P4-bmv2. (2017). <https://github.com/p4lang/behavioral-model>
- [14] P4 Language Consortium. 2017. The reference P4 program switch.p4. (2017). <https://github.com/p4lang/switch>
- [15] P4 Language Consortium. 2017. The sample P4 programs. (2017). [https://github.com/p4lang/p4c-bm/tree/master/tests/p4\\_programs](https://github.com/p4lang/p4c-bm/tree/master/tests/p4_programs)
- [16] Huynh Tu Dang, Daniele Sciascia, Marco Canini, Fernando Pedone, and Robert Soulé. 2015. NetPaxos: Consensus at Network Speed. In *Proceedings of the 1st ACM SIGCOMM Symposium on Software Defined Networking Research (SOSR '15)*. ACM, New York, NY, USA, Article 5, 7 pages. <https://doi.org/10.1145/2774993.2774999>
- [17] Nate Foster, Rob Harrison, Michael J. Freedman, Christopher Monsanto, Jennifer Rexford, Alec Story, and David Walker. 2011. Frenetic: A Network Programming Language. In *Proceedings of the 16th ACM SIGPLAN International Conference on Functional Programming (ICFP '11)*. ACM, New York, NY, USA, 279–291. <https://doi.org/10.1145/2034773.2034812>
- [18] Michael R. Garey and David S. Johnson. 1990. *Computers and Intractability; A Guide to the Theory of NP-Completeness*. W. H. Freeman & Co., New York, NY, USA.
- [19] Glen Gibb, George Varghese, Mark Horowitz, and Nick McKeown. 2013. Design Principles for Packet Parsers. In *Proceedings of the Ninth ACM/IEEE Symposium on Architectures for Networking and Communications Systems (ANCS '13)*. IEEE Press, Piscataway, NJ, USA, 13–24. <http://dl.acm.org/citation.cfm?id=2537857.2537860>
- [20] Ramesh Govindan, Ina Minei, Mahesh Kallahalla, Bikash Koley, and Amin Vahdat. 2016. Evolve or Die: High-Availability Design Principles Drawn from Google's Network Infrastructure. In *Proceedings of the 2016 ACM SIGCOMM Conference (SIGCOMM '16)*. ACM, New York, NY, USA, 58–72. <https://doi.org/10.1145/2934872.2934891>
- [21] Arpit Gupta, Laurent Vanbever, Muhammad Shahbaz, Sean P. Donovan, Brandon Schlinker, Nick Feamster, Jennifer Rexford, Scott Shenker, Russ Clark, and Ethan Katz-Bassett. 2014. SDX: A Software Defined Internet Exchange. In *Proceedings of the 2014 ACM Conference on SIGCOMM (SIGCOMM '14)*. ACM, New York, NY, USA, 551–562. <https://doi.org/10.1145/2619239.2626300>
- [22] David Hancock and Jacobus van der Merwe. 2016. HyPer4: Using P4 to Virtualize the Programmable Data Plane. In *Proceedings of the 12th International Conference on Emerging Networking Experiments and Technologies (CoNEXT '16)*. ACM, New York, NY, USA, 35–49. <https://doi.org/10.1145/2999572.2999607>
- [23] Chengchen Hu, Ji Yang, Hongbo Zhao, and Jiahua Lu. 2014. Design of All Programmable Innovation Platform for Software Defined Networking. In *Presented as part of the Open Networking Summit 2014 (ONS 2014)*. USENIX, Santa Clara, CA. <https://www.usenix.org/conference/ons2014/technical-sessions/presentation/hu-chengchen>
- [24] Brijnesh J. Jain and Klaus Obermayer. 2011. Extending Bron Kerbosch for Solving the Maximum Weight Clique Problem. *CoRR* abs/1101.1266 (2011). arXiv:1101.1266 <http://arxiv.org/abs/1101.1266>
- [25] Mikel Jimenez and Henry Kwok. 2017. Building Express Backbone: Facebook's new long-haul network. <https://code.facebook.com/posts/1782709872057497/building-express-backbone-facebook-s-new-long-haul-network/>. (2017).
- [26] Xin Jin, Jennifer Gossels, Jennifer Rexford, and David Walker. 2015. CoVisor: A Compositional Hypervisor for Software-defined Networks. In *Proceedings of the 12th USENIX Conference on Networked Systems Design and Implementation (NSDI'15)*. USENIX Association, Berkeley, CA, USA, 87–101. <http://dl.acm.org/citation.cfm?id=2789770.2789777>
- [27] Lavanya Jose, Lisa Yan, George Varghese, and Nick McKeown. 2015. Compiling Packet Programs to Reconfigurable Switches. In *Proceedings of the 12th USENIX Conference on Networked Systems Design and Implementation (NSDI'15)*. USENIX Association, Berkeley, CA, USA, 103–115. <http://dl.acm.org/citation.cfm?id=2789770.2789778>
- [28] Georgios P. Katsikas, Marcel Enguehard, Maciej Kuźniar, Gerald Q. Maguire Jr, and Dejan Kostić. 2016. SNF: synthesizing high performance NFV service chains. *PeerJ Computer Science* 2 (Nov. 2016), e98. <https://doi.org/10.7717/peerj-cs.98>
- [29] Naga Praveen Katta, Jennifer Rexford, and David Walker. 2013. Incremental Consistent Updates. In *Proceedings of the Second ACM SIGCOMM Workshop on Hot Topics in Software Defined Networking (HotSDN '13)*. ACM, New York, NY, USA, 49–54. <https://doi.org/10.1145/2491185.2491191>
- [30] Eric Keller and Evan Green. 2008. Virtualizing the Data Plane Through Source Code Merging. In *Proceedings of the ACM Workshop on Programmable Routers for Extensible Services of Tomorrow (PRESTO '08)*. ACM, New York, NY, USA, 9–14. <https://doi.org/10.1145/1397718.1397721>
- [31] Eric Keller, Minlan Yu, Matthew Caesar, and Jennifer Rexford. 2009. Virtually Eliminating Router Bugs. In *Proceedings of the 5th International Conference on Emerging Networking Experiments and Technologies (CoNEXT '09)*. ACM, New York, NY, USA, 13–24. <https://doi.org/10.1145/1658939.1658942>
- [32] Bojie Li, Zhenyuan Ruan, Wencong Xiao, Yuanwei Lu, Yongqiang Xiong, Andrew Putnam, Enhong Chen, and Lintao Zhang. 2017. KV-Direct: High-Performance In-Memory Key-Value Store with Programmable NIC. In *Proceedings of the 26th Symposium on Operating Systems Principles (SOSP '17)*. ACM, New York, NY, USA, 137–152. <https://doi.org/10.1145/3132747.3132756>
- [33] Hongqiang Harry Liu, Xin Wu, Ming Zhang, Lihua Yuan, Roger Wattenhofer, and David Maltz. 2013. zUpdate: Updating Data Center Networks with Zero Loss. In *Proceedings of the ACM SIGCOMM 2013 Conference on SIGCOMM (SIGCOMM '13)*. ACM, New York, NY, USA, 411–422. <https://doi.org/10.1145/2486001.2486005>
- [34] Shouxi Luo, Hongfang Yu, and Laurent Vanbever. 2017. Swing State: Consistent Updates for Stateful and Programmable Data Planes. In *Proceedings of the Symposium on SDN Research (SOSR '17)*. ACM, New York, NY, USA, 115–121. <https://doi.org/10.1145/3050220.3050233>
- [35] Ratul Mahajan and Roger Wattenhofer. 2013. On Consistent Updates in Software Defined Networks. In *Proceedings of the Twelfth ACM Workshop on Hot Topics in Networks (HotNets-XII)*. ACM, New York, NY, USA, Article 20, 7 pages. <https://doi.org/10.1145/2535771.2535791>
- [36] Rui Miao, Hongyi Zeng, Changhoon Kim, Jeongkeun Lee, and Minlan Yu. 2017. SilkRoad: Making Stateful Layer-4 Load Balancing Fast and Cheap Using Switching ASICs. In *Proceedings of the Conference of the ACM Special Interest Group on Data Communication (SIGCOMM '17)*. ACM, New York, NY, USA, 15–28. <https://doi.org/10.1145/3098822.3098824>
- [37] Tal Mizrahi, Efi Saat, and Yoram Moses. 2015. Timed Consistent Network Updates. In *Proceedings of the 1st ACM SIGCOMM Symposium on Software Defined Networking Research (SOSR '15)*. ACM, New York, NY, USA, Article 21, 14 pages. <https://doi.org/10.1145/2774993.2775001>
- [38] Christopher Monsanto, Joshua Reich, Nate Foster, Jennifer Rexford, and David Walker. 2013. Composing Software Defined Networks. In *10th USENIX Symposium on Networked Systems Design and Implementation (NSDI 13)*. USENIX Association, Lombard, IL, 1–13. <https://www.usenix.org/conference/nsdi13/technical-sessions/presentation/monsanto>
- [39] Barefoot Networks. 2016. Barefoot Whitepaper: The World's Fastest and Most Programmable Networks. (2016). <https://barefootnetworks.com/resources/worlds-fastest-most-programmable-networks/>
- [40] Bruno Nogueira, Rian G. S. Pinheiro, and Anand Subramanian. 2018. A hybrid iterated local search heuristic for the maximum weight independent set problem. *Optimization Letters* 12, 3 (01 May 2018). <https://doi.org/10.1007/s11590-017-1128-7>
- [41] Recep Ozdag. 2012. Intel® Ethernet Switch FM6000 Series-Software Defined Networking. (2012).
- [42] D. L. Parnas. 1972. On the Criteria to Be Used in Decomposing Systems into Modules. *Commun. ACM* 15, 12 (Dec. 1972), 1053–1058. <https://doi.org/10.1145/361598.361623>
- [43] Dhrubajyoti Saha, Abhishek Samanta, and Smruti R Sarangi. 2009. Theoretical framework for eliminating redundancy in workflows. In *Services Computing, 2009. SCC '09. IEEE International Conference on*. IEEE, 41–48.
- [44] Shuichi Sakai, Mitsunori Togasaki, and Koichi Yamazaki. 2003. A note on greedy algorithms for the maximum weighted independent set problem. *Discrete Applied Mathematics* 126, 2 (2003), 313 – 322. [https://doi.org/10.1016/S0166-218X\(02\)00205-6](https://doi.org/10.1016/S0166-218X(02)00205-6)
- [45] Rob Sherwood, Glen Gibb, Kok-Kiong Yap, Guido Appenzeller, Martin Casado, Nick McKeown, and Guru Parulkar. 2010. Can the Production Network Be the Testbed?. In *Proceedings of the 9th USENIX Conference on Operating Systems Design and Implementation (OSDI'10)*. USENIX Association, Berkeley, CA, USA, 365–378. <http://dl.acm.org/citation.cfm?id=1924943.1924969>
- [46] Arjun Singh, Joon Ong, Amit Agarwal, Glen Anderson, Ashby Armistead, Roy Bannon, Seb Boving, Gaurav Desai, Bob Felderman, Paulie Germano, Anand Kana-gala, Jeff Provost, Jason Simmons, Eiichi Tada, Jim Wanderer, Urs Hölzle, Stephen

- Stuart, and Amin Vahdat. 2015. Jupiter Rising: A Decade of Clos Topologies and Centralized Control in Google's Datacenter Network. In *Proceedings of the 2015 ACM Conference on Special Interest Group on Data Communication (SIGCOMM '15)*. ACM, New York, NY, USA, 183–197. <https://doi.org/10.1145/2785956.2787508>
- [47] Anirudh Sivaraman, Alvin Cheung, Mihai Budiu, Changhoon Kim, Mohammad Alizadeh, Hari Balakrishnan, George Varghese, Nick McKeown, and Steve Licking. 2016. Packet Transactions: High-Level Programming for Line-Rate Switches. In *Proceedings of the 2016 ACM SIGCOMM Conference (SIGCOMM '16)*. ACM, New York, NY, USA, 15–28. <https://doi.org/10.1145/2934872.2934900>
- [48] Vibhaalakshmi Sivaraman, Srinivas Narayana, Ori Rottenstreich, S. Muthukrishnan, and Jennifer Rexford. 2017. Heavy-Hitter Detection Entirely in the Data Plane. In *Proceedings of the Symposium on SDN Research (SOSR '17)*. ACM, New York, NY, USA, 164–176. <https://doi.org/10.1145/3050220.3063772>
- [49] Hardik Soni, Thierry Turlletti, and Walid Dabbous. 2018. P4Bricks: Enabling multiprocessing using Linker-based network data plane architecture. (Feb. 2018). <https://hal.inria.fr/hal-01632431> working paper.
- [50] Peng Sun, Ratul Mahajan, Jennifer Rexford, Lihua Yuan, Ming Zhang, and Ahsan Arefin. 2014. A Network-state Management Service. *SIGCOMM Comput. Commun. Rev.* 44, 4 (Aug. 2014), 563–574. <https://doi.org/10.1145/2740070.2626298>
- [51] Yu-Wei Eric Sung, Xiaozheng Tie, Starsky H.Y. Wong, and Hongyi Zeng. 2016. Robotron: Top-down Network Management at Facebook Scale. In *Proceedings of the 2016 ACM SIGCOMM Conference (SIGCOMM '16)*. ACM, New York, NY, USA, 426–439. <https://doi.org/10.1145/2934872.2934874>
- [52] Kaushik Veeraraghavan, Justin Meza, David Chou, Wonho Kim, Sonia Margulis, Scott Michelson, Rajesh Nishtala, Daniel Obenshain, Dmitri Perelman, and Yee Jiun Song. 2016. Kraken: Leveraging Live Traffic Tests to Identify and Resolve Resource Utilization Bottlenecks in Large Scale Web Services. In *12th USENIX Symposium on Operating Systems Design and Implementation (OSDI 16)*. USENIX Association, Savannah, GA, 635–651. <https://www.usenix.org/conference/osdi16/technical-sessions/presentation/veeraraghavan>
- [53] Han Wang, Robert Soulé, Huynh Tu Dang, Ki Suh Lee, Vishal Shrivastav, Nate Foster, and Hakim Weatherspoon. 2017. P4FPGA: A Rapid Prototyping Framework for P4. In *Proceedings of the Symposium on SDN Research (SOSR '17)*. ACM, New York, NY, USA, 122–135.
- [54] Xilinx. 2014. SDNet. (2014). <http://www.xilinx.com/products/design-tools/software-zone/sdnet.html>
- [55] Xilinx. 2017. Ternary Content Addressable Memory (TCAM) Search IP for SDNet SmartCORE IP Product Guide. (2017).
- [56] Cheng Zhang, Jun Bi, Yu Zhou, Abdul Basit Dogar, and Jianping Wu. 2017. HyperV: A High Performance Hypervisor for Virtualization of the Programmable Data Plane. In *2017 26th International Conference on Computer Communication and Networks (ICCCN)*, 1–9. <https://doi.org/10.1109/ICCCN.2017.8038396>
- [57] Cheng Zhang, Jun Bi, Yu Zhou, Abdul Basit Dogar, and Jianping Wu. 2017. MPVisor: A Modular Programmable Data Plane Hypervisor. In *Proceedings of the Symposium on SDN Research (SOSR '17)*. ACM, New York, NY, USA, 179–180. <https://doi.org/10.1145/3050220.3060600>
- [58] Ying Zhang, Neda Beheshti, and Ravi Manghirmalani. 2014. NetRevert: Rollback Recovery in SDN. In *Proceedings of the Third Workshop on Hot Topics in Software Defined Networking (HotSDN '14)*. ACM, New York, NY, USA, 231–232. <https://doi.org/10.1145/2620728.2620779>
- [59] Danyang Zhuo, Qiao Zhang, Xin Yang, and Vincent Liu. 2016. Canaries in the Network. In *Proceedings of the 15th ACM Workshop on Hot Topics in Networks (HotNets '16)*. ACM, New York, NY, USA, 36–42. <https://doi.org/10.1145/3005745.3005767>